

Estimation of causal orders in a linear non-Gaussian acyclic model: a method robust against latent confounders

Tatsuya Tashiro¹, Shohei Shimizu¹, Aapo Hyvärinen², and Takashi Washio¹

¹ The Institute of Scientific and Industrial Research, Osaka University, Japan

² Dept. of Mathematics and Statistics, Dept. of Computer Science / HIIT, University of Helsinki, Finland

Abstract. We consider to learn a causal ordering of variables in a linear non-Gaussian acyclic model called LiNGAM. Several existing methods have been shown to consistently estimate a causal ordering assuming that all the model assumptions are correct. But, the estimation results could be distorted if some assumptions actually are violated. In this paper, we propose a new algorithm for learning causal orders that is robust against one typical violation of the model assumptions: latent confounders. We demonstrate the effectiveness of our method using artificial data.

Keywords: Bayesian networks, causal discovery, non-Gaussianity, latent confounders, independent component analysis

1 Introduction

Bayesian networks have been widely used to analyze causal relations of variables in many empirical sciences [11]. A common assumption is linear-Gaussianity. But this poses serious identifiability problems so that many important models are indistinguishable with no prior knowledge on the structures. Recently, it was shown [9] that use of non-Gaussianity allows the full structure of a linear acyclic model to be identified without pre-specifying any causal orders of variables. The new model, a Linear Non-Gaussian Acyclic Model called LiNGAM [9], is closely related to independent component analysis (ICA) [7].

Existing estimation methods [9,10] for LiNGAM learn causal orders assuming that all the model assumptions hold. Therefore, these algorithms could return completely wrong estimation results when some of the model assumptions is violated. Thus, in this paper, we propose a new algorithm for learning causal orders that is robust against one typical model violation, *i.e.*, latent confounders. A latent confounder means a variable which is not observed but which exerts a causal influence on some of the observed variables.

The paper is organized as follows. We first review LiNGAM [9] and its extension to latent confounder cases [6] in Section 2. In Section 3, we propose a new algorithm to learn causal orders in LiNGAM with latent confounders. Simulations are conducted in Section 4. We conclude this paper in Section 5.

2 Background: LiNGAM with latent confounders

We briefly review a linear non-Gaussian acyclic model called LiNGAM [9] and an extension of the LiNGAM to cases with latent confounding variables [6].

In LiNGAM [9], causal relations of observed variables x_i are modeled as:

$$x_i = \sum_{k(j) < k(i)} b_{ij} x_j + e_i, \quad (1)$$

where $k(i)$ is such a causal ordering of variables x_i that they graphically form a directed acyclic graph (DAG) so that no later variable determines, *i.e.*, has a directed path on any earlier variable, e_i are external influences, and b_{ij} are connection strengths. In matrix form, the model (1) is written as

$$\mathbf{x} = \mathbf{B}\mathbf{x} + \mathbf{e}, \quad (2)$$

where the connection strength matrix \mathbf{B} collects b_{ij} and the vectors \mathbf{x} and \mathbf{e} collect x_i and e_i . Note that the matrix \mathbf{B} can be permuted to be lower triangular with all zeros on the diagonal if simultaneous equal row and column permutations are made according to a causal ordering $k(i)$ due to the acyclicity. The zero/non-zero pattern of b_{ij} corresponds to the absence/existence pattern of directed edges. External influences e_i follow non-Gaussian continuous distributions with zero mean and non-zero variance and are mutually independent. The non-Gaussianity assumption on e_i enables identification of a causal ordering $k(i)$ based on data \mathbf{x} only [9]. This feature is a big advantage over conventional Bayesian networks based on the Gaussianity assumption on e_i [11].

Next, LiNGAM with latent confounders [6] can be formulated as follows:

$$\mathbf{x} = \mathbf{B}\mathbf{x} + \mathbf{\Lambda}\mathbf{f} + \mathbf{e}, \quad (3)$$

where the difference with LiNGAM (2) is the existence of latent confounding variable vector \mathbf{f} . A latent confounding variable is such an latent variable that is a parent of more than or equal to two observed variables. The vector \mathbf{f} collects non-Gaussian latent confounders f_j with zero mean and non-zero variance ($j = 1, \dots, q$). Without loss of generality [6], latent confounders f_j are assumed to be mutually independent. The matrix $\mathbf{\Lambda}$ collects λ_{ij} which denotes the connection strength from f_j to x_i . For each j , at least two λ_{ij} are non-zero since a latent confounder is defined to have at least two children. Further, it is assumed [6] that correlation and conditional correlation of x_i , f_i and e_i are entailed by the graph structure only, *i.e.*, the zero/non-zero status of b_{ij} and λ_{ij} . This is a well-known assumption called faithfulness in causal discovery [11].

The central problem of causal discovery based on the latent variable LiNGAM in Equation (3) is to estimate *as many* of causal orders $k(i)$ and connection strengths b_{ij} *as possible* based on data \mathbf{x} only. This is because in many cases only an equivalence class of the true model whose members produce the exact same observed distribution is identifiable [6].

In [6], an estimation method based on overcomplete ICA was proposed. However, overcomplete ICA methods are often not very reliable and get stuck in local

optima. Thus, in [2], a method that does not use overcomplete ICA was proposed to first find variable *pairs* that are not affected by latent confounders and then estimate a causal ordering of one to the other instead of a causal ordering of more than two variables.

3 A hybrid estimation approach

In this section, we propose a new approach for estimating causal orders of more than two variables without explicitly modeling latent confounders. We first provide principles to identify such an exogenous (root) variable and a sink variable that are not affected by latent confounders in the latent variable LiNGAM in Equation (3) (if such variables exist) and next present an estimation algorithm. Recent estimation methods [8, 10] for LiNGAM in Equation (2) and its non-linear extension [5] learn a causal ordering by finding causal orders one by one either from the top downward or from the bottom upward assuming no latent confounders. We extend these ideas to latent confounder cases.

We first generalize Lemma 1 of [10] for the case of latent confounders.

Lemma 1 *Assume that all the model assumptions of the latent variable LiNGAM (3) are met and the sample size is infinite. Denote by $r_i^{(j)}$ the residuals when x_i are regressed on x_j : $r_i^{(j)} = x_i - \frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)} x_j$ ($i \neq j$). Then a variable x_j is an exogenous variable in the sense that it has no parent observed variable nor latent confounder if and only if x_j is independent of its residuals $r_i^{(j)}$ for all $i \neq j$. \square*

Next, we generalize the idea of [8] for the case of latent confounders.

Lemma 2 *Assume that all the model assumptions of the latent variable LiNGAM (3) are met and the sample size is infinite. Denote by $\mathbf{x}_{(-j)}$ a vector that contains all the variables other than x_j . Denote by $r_j^{(-j)}$ the residual when x_j is regressed on $\mathbf{x}_{(-j)}$, i.e., $r_j^{(-j)} = x_j - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \mathbf{x}_{(-j)}$, where $\Sigma = \begin{bmatrix} \sigma_j & \boldsymbol{\sigma}_{j(-j)}^T \\ \boldsymbol{\sigma}_{j(-j)} & \Sigma_{(-j)} \end{bmatrix}$ is the covariance matrix of $[x_j, \mathbf{x}_{(-j)}^T]^T$. Then a variable x_j is a sink variable in the sense that it has no child observed variable nor latent confounder if and only if $\mathbf{x}_{(-j)}$ is independent of its residual $r_j^{(-j)}$. \square*

The proofs of these lemmas are given in the appendix.

Thus, we can take a hybrid estimation approach that uses these two principles. We first identify an exogenous variable by finding a variable that is most independent of its residuals and remove the effect of the exogenous variable from the other variables by regressing it out. We repeat this until independence of every variable and any of its residuals is statistically rejected. Dependency between every variable and any of its residuals implies that such an exogenous variable in Lemma 1 does not exist or some model assumption of latent variable LiNGAM (3) is violated. Similarly, we next identify a sink variable in the remaining variables by finding a variable that its regressors and its residual are

most independent and disregard the sink variable. We repeat this until independence is statistically rejected for every variable. We test pairwise independence between variables and the residuals using a kernel-based independence measure called HSIC [4] and combine the resulting p -values using a well-known Fisher's method [3]. We use Bonferroni correction for multiple comparison dividing the significance level by the maximum number of tests $p-1$.

Thus, the estimation consists of the following steps:

1. Given a p -dimensional random vector \mathbf{x} , a set of its variable subscripts U , a $p \times n$ data matrix of the random vector as \mathbf{X} and a significance level α , initialize an ordered list of variables $K_{head} := \emptyset$ and $K_{tail} := \emptyset$ and $m := 1$. K_{head} and K_{tail} denote first $|K_{head}|$ orders of variables and last $|K_{tail}|$ orders of variables respectively, where each of $|K_{head}|$ and $|K_{tail}|$ denotes the number of elements in the list.
2. Let $\tilde{\mathbf{x}} = \mathbf{x}$ and $\tilde{\mathbf{X}} = \mathbf{X}$ and find causal orders one by one from the top downward:

- (a) Do the following steps for all $j \in U \setminus K_{head}$: Perform least squares regressions of \tilde{x}_i on \tilde{x}_j for all $i \in U \setminus K_{head}$ ($i \neq j$) and compute the residual vectors $\tilde{\mathbf{r}}^{(j)}$. Then, find a variable \tilde{x}_m that is most independent of its residuals:

$$\tilde{x}_m = \arg \max_{j \in U \setminus K_{head}} P_{Fisher}(\tilde{x}_j, \tilde{\mathbf{r}}^{(j)}), \quad (4)$$

where $P_{Fisher}(\tilde{x}_j, \tilde{\mathbf{r}}^{(j)})$ is the p -value of the test statistic defined as $-\frac{2}{\sum_i \log\{P_H(\tilde{x}_j, \tilde{r}_i^{(j)})\}}$, where $P_H(\tilde{x}_j, \tilde{r}_i^{(j)})$ is the p -value of the HSIC.

- (b) Go to Step 3 if $P_{Fisher}(\tilde{x}_m, \tilde{\mathbf{r}}^{(m)}) < \alpha/(p-1)$.
- (c) Append m to the end of K_{head} and let $\tilde{\mathbf{x}} := \tilde{\mathbf{r}}^{(m)}$ and $\tilde{\mathbf{X}} := \tilde{\mathbf{R}}^{(m)}$. If $|K_{head}| = p-1$, append the remaining variable subscript to the end of K_{head} and go to Step 4. Otherwise, go back to Step (2a).
3. If $|K_{head}| < p-2$,³ let $\mathbf{x}' = \mathbf{x}$ and $\mathbf{X}' = \mathbf{X}$ and $U' := U \setminus K_{head}$ and find causal orders one by one from the bottom upward:
 - (a) Do the following steps for all $j \in U' \setminus K_{tail}$: Collect all the variables except x'_j in a vector $\mathbf{x}'_{(-j)}$. Perform least squares regressions of x'_j on $\mathbf{x}'_{(-j)}$ and compute the residual $r'^{(-j)}_j$. Then, find such a variable x'_m that its regressors and its residual are most independent:

$$x'_m = \arg \max_{j \in U' \setminus K_{tail}} P_{Fisher}(\mathbf{x}'_{(-j)}, r'^{(-j)}_j). \quad (5)$$

- (b) Go to Step 4 if $P_{Fisher}(\mathbf{x}'_{(-m)}, r'^{(-m)}_m) < \alpha/(p-1)$.
- (c) Append m to the top of K_{tail} and let $\mathbf{x}' = \mathbf{x}'_{(-m)}$ and $\mathbf{X}' = \mathbf{X}'_{(-m)}$. Go to Step 4 if $|U' \setminus K_{tail}| < 3$.³ Otherwise go back to Step (3a).
4. Estimate connection strengths b_{ij} for variables in K_{head} and K_{tail} by doing multiple regression of every variable x_i in K_{head} and K_{tail} on all of its non-descendants x_j with $k(j) < k(i)$.

³ We do not examine remaining two variables in Step 3 since it is already implied in Step 2 that some latent confounders exist or some model assumption is violated.

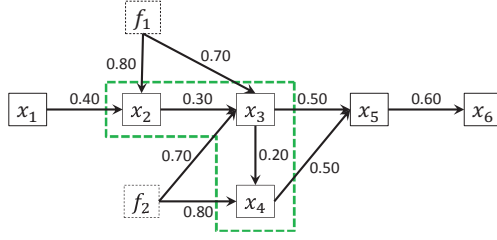


Fig. 1. True network used in the simulation. The variables f_1 and f_2 were latent confounders. The green contours include variables that share f_1 or f_2 . The external influences e_i are omitted to be shown.

Note that our algorithm would output no causal orders in cases that such exogenous variables and sink variables as in Lemmas 1 and 2 do not exist, although the outputs are still correct. One way to learn more causal orders in those cases would be to develop a divide-and-conquer algorithm that divides variables into subsets where such exogenous or sink variables exist and integrates the estimation results on the subsets. This is an important direction of future research.

4 Experiments on artificial data

We compared our method with an estimation method for LiNGAM (2) called DirectLiNGAM [10] that does not allow latent confounders and an estimation method for latent variable LiNGAM (3) called Pairwise LvLiNGAM [2]. If there is no latent confounders, all the methods should estimate correct causal orders for enough large sample sizes. The number of variables was 6, and the sample sizes tested were 500, 1000, 2000. The original network used was shown in Figure 1. The e_1, e_4 and f_1 followed a multimodal asymmetric mixture of two Gaussians, e_2, e_5, f_2 followed a double exponential distribution, and e_3 and e_6 followed a multimodal symmetric mixture of two Gaussians. The standard deviations of the e_i were set so that their signal-to-noise ratios, *i.e.*, $\text{var}(x_i)/\text{var}(e_i) - 1$ were all ones. The number of trials was 100. The significance level α was 0.05.

First, to evaluate performance of estimating causal orders $k(i)$, we computed the percentage of correctly estimated causal orders in estimated causal orders between two variables (Precision) and the percentage of correctly estimated causal orders in actual causal orders between two variables that share no latent confounders in the true data generating network (Recall). The reason why only pairwise causal orders were evaluated was that Pairwise LvLiNGAM only estimates causal orders of two variables unlike our method and DirectLiNGAM. Tables 1 and 2 show the results. Regarding precisions, our method was comparable to Pairwise LvLiNGAM and the two methods were much better than DirectLiNGAM for all the conditions. Regarding recalls, our method was better than both DirectLiNGAM and Pairwise LvLiNGAM for all the conditions.

Table 1. Precisions

	Sample size		
	500	1000	2000
Our method	0.78	0.80	0.80
DirectLiNGAM	0.65	0.64	0.64
Pairwise LvLiNGAM	0.79	0.81	0.81

Table 2. Recalls

	Sample size		
	500	1000	2000
Our method	0.97	0.99	0.99
DirectLiNGAM	0.81	0.80	0.81
Pairwise LvLiNGAM	0.86	0.89	0.90

Next, to evaluate the performance in estimating connection strengths b_{ij} , we computed the root mean square errors between true connection strengths and estimated ones. The root mean square errors for our method and DirectLiNGAM were 0.079 and 0.090 for 500 data points, 0.070 and 0.079 for 1000 data points and 0.015 and 0.057 for 2000 data points, respectively, where our method was more accurate. Note that Pairwise LvLiNGAM does not estimate b_{ij} .

5 Conclusions

We proposed a new algorithm for learning causal orders, which is robust against latent confounders. In experiments on artificial data, our approach learned more causal orders accurately than two existing methods. In future work, we would like to test our method on real-world data including functional magnetic resonance imaging data to analyze causal interactions between brain regions.

Acknowledgments. S.S and T.W. were supported by KAKENHI #24700275 and #22300054. We thank Patrik Hoyer and Doris Entner for helpful comments.

References

1. Darmonis, G.: Analyse générale des liaisons stochastiques. Review of the International Statistical Institute 21, 2–8 (1953)
2. Entner, D., Hoyer, P.O.: Discovering unconfounded causal relationships using linear non-gaussian models. In: New Frontiers in Artificial Intelligence, Lecture Notes in Computer Science. vol. 6797, pp. 181–195 (2011)
3. Fisher, R.: Statistical methods for research workers. Oliver and Boyd (1950)
4. Gretton, A., Fukumizu, K., Teo, C., Song, L., Schölkopf, B., Smola, A.J.: A kernel statistical test of independence. In: Advances in Neural Information Processing Systems 20. MIT Press, Cambridge, MA (2008)
5. Hoyer, P.O., Janzing, D., Mooij, J., Peters, J., Schölkopf, B.: Nonlinear causal discovery with additive noise models. In: Advances in Neural Information Processing Systems 21, pp. 689–696 (2009)
6. Hoyer, P.O., Shimizu, S., Kerminen, A., Palviainen, M.: Estimation of causal effects using linear non-gaussian causal models with hidden variables. International Journal of Approximate Reasoning 49(2), 362–378 (2008)
7. Hyvärinen, A., Karhunen, J., Oja, E.: Independent component analysis. Wiley, New York (2001)

8. Mooij, J., Janzing, D., Peters, J., Schölkopf, B.: Regression by dependence minimization and its application to causal inference in additive noise models. In: Proc. the 26th Int. Conf. on Machine Learning (ICML2009). pp. 745–752 (2009)
9. Shimizu, S., Hoyer, P.O., Hyvärinen, A., Kerminen, A.: A linear non-gaussian acyclic model for causal discovery. J. Mach. Learn. Res. 7, 2003–2030 (2006)
10. Shimizu, S., Inazumi, T., Sogawa, Y., Hyvärinen, A., Kawahara, Y., Washio, T., Hoyer, P.O., Bollen, K.: DirectLiNGAM: A direct method for learning a linear non-Gaussian structural equation model. J. Mach. Learn. Res. 12, 1225–1248 (2011)
11. Spirtes, P., Glymour, C., Scheines, R.: Causation, Prediction, and Search. Springer Verlag (1993), (2nd ed. MIT Press 2000)

Appendix: Proofs of the lemmas

Theorem 1 (Darmois-Skitovitch theorem (D-S theorem) [1]) *Define two random variables y_1 and y_2 as linear combinations of independent random variables $s_i (i=1, \dots, q)$: $y_1 = \sum_{i=1}^q \alpha_i s_i$, $y_2 = \sum_{i=1}^q \beta_i s_i$. Then, if y_1 and y_2 are independent, all variables s_j for which $\alpha_j \beta_j \neq 0$ are Gaussian.* \square

In other words, this theorem means that if there exists a non-Gaussian s_j for which $\alpha_j \beta_j \neq 0$, y_1 and y_2 are dependent.

Further, Lemma 3 of [2] has shown that the regressor and its residual in simple linear regression are dependent if there are some latent confounders between the regressor and regressand in the latent variable LiNGAM (3).

Proof of Lemma 1 i) Assume that x_j has at least one parent observed variable or latent confounder. Let P_j denote the set of the parent variables of x_j . Then one can write $x_j = \sum_{p_h \in P_j} w_{jh} p_h + e_j$, where the parent variables p_h are independent of e_j and the coefficients w_{jh} are non-zero. Suppose that x_i is a parent of x_j . For such x_i , we have $r_i^{(j)} = x_i - \frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)} x_j = x_i - \frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)} (\sum_{p_h \in P_j} w_{jh} p_h + e_j) = \left\{ 1 - \frac{w_{ji} \text{cov}(x_i, x_j)}{\text{var}(x_j)} \right\} x_i - \frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)} \sum_{p_h \in P_j, p_h \neq x_i} w_{jh} p_h - \frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)} e_j$. Each of those parent variables (including x_i) in P_j is a linear combination of external influences *other than* e_j and latent confounders that are non-Gaussian and independent. Thus, the $r_i^{(j)}$ and x_j can be written as linear combinations of non-Gaussian and independent external influences including e_j and latent confounders. Further, the coefficient of e_j on $r_i^{(j)}$ is non-zero since $\text{cov}(x_i, x_j) \neq 0$ due to the faithfulness and that on x_j is one by definition. These imply that $r_i^{(j)}$ and x_j are dependent since $r_i^{(j)}$, x_j and e_j correspond to y_1 , y_2 , s_j in D-S theorem, respectively. Next, for the other case that x_j has a latent confounder, $r_i^{(-j)}$ and an observed variable can be shown to be dependent using Lemma 3 of [2] since by definition at least one observed variable shares the latent confounder with x_j .

ii) The converse of contrapositive of i) is straightforward using the model definition. From i) and ii), the lemma is proven. \blacksquare

Proof of Lemma 2 i) Assume that a variable x_j has at least one child observed variable or latent confounder. First, without loss of generality, one can write

$$\mathbf{x} = \begin{bmatrix} x_j \\ \mathbf{x}_{(-j)} \end{bmatrix} = (\mathbf{I} - \mathbf{B})^{-1}(\mathbf{A}\mathbf{f} + \mathbf{e}) = \mathbf{A}(\mathbf{A}\mathbf{f} + \mathbf{e}) \quad (6)$$

$$= \begin{bmatrix} 1 & \mathbf{a}_{j(-j)}^T \\ \mathbf{a}_{(-j)j} & \mathbf{A}_{(-j)} \end{bmatrix} \begin{bmatrix} \boldsymbol{\lambda}_j^T \mathbf{f} + e_j \\ \mathbf{A}_{(-j)} \mathbf{f} + \mathbf{e}_{(-j)} \end{bmatrix}, \quad (7)$$

where each of \mathbf{A} ($= (\mathbf{I} - \mathbf{B})^{-1}$) and $\mathbf{A}_{(-j)}$ is invertible and can be permuted to be a lower triangular matrix with the diagonal elements being ones if the rows and columns are simultaneously permuted according to the causal ordering $k(i)$. The same applies to the inverse of \mathbf{A} :

$$\mathbf{A}^{-1} = \begin{bmatrix} (1 - \mathbf{a}_{j(-j)}^T \mathbf{A}_{(-j)}^{-1} \mathbf{a}_{(-j)j})^{-1} & -\mathbf{a}_{j(-j)}^T \mathbf{D}^{-1} \\ -\mathbf{D}^{-1} \mathbf{a}_{(-j)j} & \mathbf{D}^{-1} \end{bmatrix}, \quad (8)$$

where $\mathbf{D} = \mathbf{A}_{(-j)} - \mathbf{a}_{(-j)j} \mathbf{a}_{j(-j)}^T$. Thus, $1 - \mathbf{a}_{j(-j)}^T \mathbf{A}_{(-j)}^{-1} \mathbf{a}_{(-j)j} = 1$. Then,

$$r_j^{(-j)} = x_j - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \mathbf{x}_{(-j)} \quad (9)$$

$$= \boldsymbol{\lambda}_j^T \mathbf{f} + e_j + \mathbf{a}_{j(-j)}^T (\mathbf{A}_{(-j)} \mathbf{f} + \mathbf{e}_{(-j)}) - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \{ \mathbf{a}_{(-j)j} (\boldsymbol{\lambda}_j^T \mathbf{f} + e_j) + \mathbf{A}_{(-j)} (\mathbf{A}_{(-j)} \mathbf{f} + \mathbf{e}_{(-j)}) \} \quad (10)$$

$$= \{ \boldsymbol{\lambda}_j^T + \mathbf{a}_{j(-j)}^T \mathbf{A}_{(-j)} - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} (\mathbf{a}_{(-j)j} \boldsymbol{\lambda}_j^T + \mathbf{A}_{(-j)} \mathbf{A}_{(-j)}) \} \mathbf{f} + \{ 1 - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \mathbf{a}_{(-j)j} \} e_j + \{ \mathbf{a}_{j(-j)}^T - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \mathbf{A}_{(-j)} \} \mathbf{e}_{(-j)}. \quad (11)$$

In Equation (11), if $\mathbf{a}_{j(-j)}^T - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \mathbf{A}_{(-j)} = \mathbf{0}^T$, then we have

$$r_j^{(-j)} = \{ \boldsymbol{\lambda}_j^T (1 - \mathbf{a}_{j(-j)}^T \mathbf{A}_{(-j)}^{-1} \mathbf{a}_{(-j)j}) \} \mathbf{f} + \{ 1 - \mathbf{a}_{j(-j)}^T \mathbf{A}_{(-j)}^{-1} \mathbf{a}_{(-j)j} \} e_j \quad (12)$$

$$= \boldsymbol{\lambda}_j^T \mathbf{f} + e_j. \quad (13)$$

Thus, the coefficient of e_j on $r_j^{(-j)}$ is one. Now, suppose that x_j has a child x_i . The coefficient of e_j on x_i is non-zero due to the faithfulness. Thus, $r_j^{(-j)}$ and x_i are dependent due to D-S theorem. Next, suppose that x_j has a latent confounder f_i . Then, in Equation (11), the corresponding element in $\boldsymbol{\lambda}_j$ is not zero, *i.e.*, the coefficient of f_i on $r_j^{(-j)}$ is not zero. Further, f_i has a non-zero coefficient on at least one variable in $\mathbf{x}_{(-j)}$ due to the definition of latent confounders and faithfulness. Therefore, $r_j^{(-j)}$ and $\mathbf{x}_{(-j)}$ are dependent due to D-S theorem.

On the other hand, in Equation (11), if $\mathbf{a}_{j(-j)}^T - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \mathbf{A}_{(-j)} \neq \mathbf{0}^T$, at least one of the coefficients of the elements in $\mathbf{e}_{(-j)}$ on $r_j^{(-j)}$ is not zero. By definition, every element in $\mathbf{e}_{(-j)}$ has a non-zero coefficient on the corresponding element in $\mathbf{x}_{(-j)}$. Thus, $r_j^{(-j)}$ and $\mathbf{x}_{(-j)}$ are dependent due to D-S theorem.

ii) The converse of contrapositive of i) is straightforward using the model definition. From i) and ii), the lemma is proven. ■